

# DYSARTHRIC VOCAL INTERFACES WITH MINIMAL TRAINING DATA

*Jort F. Gemmeke<sup>1</sup>, Siddharth Sehgal<sup>2</sup>, Stuart Cunningham<sup>2</sup>, Hugo Van hamme<sup>1</sup>*

<sup>1</sup>ESAT-PSI, KU Leuven, Kasteelpark Arenberg 10, 3001, Leuven, Belgium

<sup>2</sup>Department of Human Communication Sciences, University of Sheffield, Sheffield, United Kingdom

email: jgemmeke@amadana.nl

## ABSTRACT

Over the past decade, several speech-based electronic assistive technologies (EATs) have been developed that target users with dysarthric speech. These EATs include vocal command & control systems, but also voice-input voice-output communication aids (VIVOCAs). In these systems, the vocal interfaces are based on automatic speech recognition systems (ASR), but this approach requires much training data and detailed annotation. In this work we evaluate an alternative approach, which works by mining utterance-based representations of speech for recurrent acoustic patterns, with the goal of achieving usable recognition accuracies with less speaker-specific training data. Comparisons with a conventional ASR system on dysarthric speech databases show that the proposed approach offers a substantial reduction in the amount of training data needed to achieve the same recognition accuracies.

**Index Terms:** vocal user interface, dysarthric speech, non-negative matrix factorisation

## 1. INTRODUCTION

Spoken language communication is central to daily life, but as many as 1.3% of the population cannot use natural speech to communicate reliably [1]. Impaired speech can often be unintelligible to unfamiliar communication partners, and it also can make the use of conventional voice controlled command & control (C&C) systems problematic. Such systems, however, can significantly contribute to the independence of living and quality of life of users with restricted motor control [2].

Over the past decade, several speech-based electronic assistive technologies (EATs) have been developed that target users with dysarthric speech. These EATs include vocal C&C systems [3, 4], but also voice-input voice-output communication aids (VIVOCAs) [5]. The three challenges these systems face are that 1) The number of phones that can be produced is often severely restricted, making it difficult to distinguish between words, 2) dysarthric speech varies greatly between speakers and 3) speaking often requires great effort, thus restricting the amount of training or adaption material that can be collected.

Conventional EATs for dysarthric speech are based on automatic speech recognition (ASR), employing either speaker-independent

acoustic models trained on a large corpus with adaptation to the target speaker [6, 7, 8, 9], or speaker-dependent models trained directly on speech material from the target user [3, 5]. Although adaptation approaches typically require less speech material from the target user than speaker-dependent modelling approaches, their performance largely depends on the exact speech characteristics.

The amount of training data required for conventional Hidden Markov Model (HMM)-based ASR is high, however; especially for speakers for whom speaking takes great effort, data collection can take many months. In this work we evaluate an alternative approach, to achieve two goals: 1) Achieving usable recognition accuracies with less training data, in order to minimize the initial effort of the target user, and 2) Achieving usable recognition accuracies with less detailed annotation - training a vocal interface using an unordered list of keywords that are contained in the sentences, rather than a word-by-word transcript. Ideally, meeting these goals results in vocal interfaces which adapt completely on-line and avoid any prior speaker-dependent data-collection. In this work, we investigate to what extent this approach can be used to augment, or even replace, a conventional speaker-dependent ASR system for dysarthric speakers.

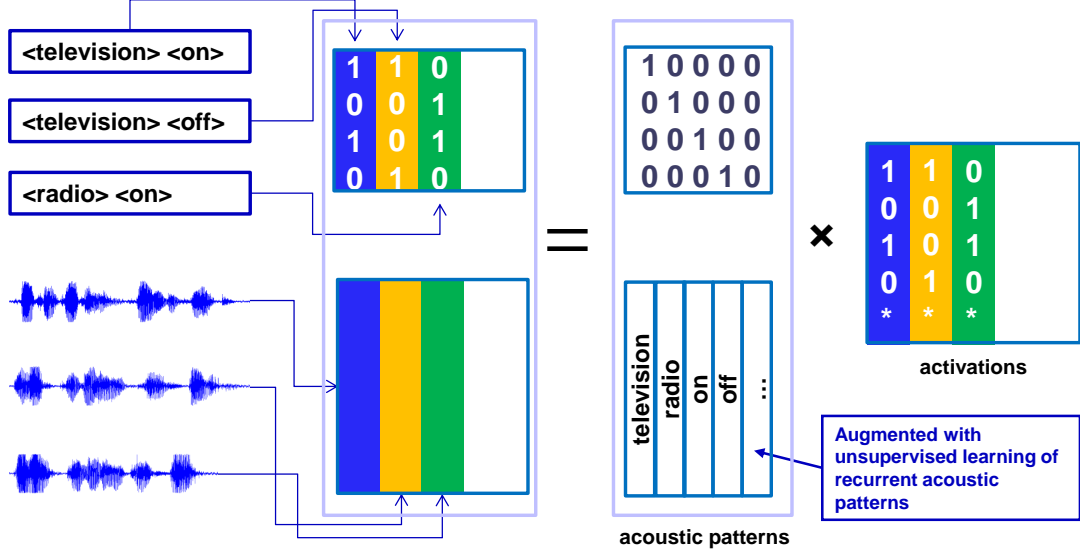
The method works by mining utterance-based representations of speech for recurrent acoustic patterns [4]. This speaker-dependent approach, developed in the ALADIN project, maps these acoustic patterns directly to (parts of) commands, which means it is language independent and does not require a pre-defined vocabulary, grammar or even knowledge of word order in the training data. The ALADIN system has been shown to yield relatively high recognition accuracies even after a single training sample of each word or command [10, 11], but the evaluations on dysarthric speech have been extremely limited and its performance has not been compared with conventional systems.

The contributions of the paper are twofold. First, we evaluate both a speaker-dependent ASR approach and the ALADIN approach on large dysarthric speech databases, with speech of severely impaired speakers, and characterise performance as a function of the amount of training data needed. Second, we evaluated the performance on both isolated words and on C&C sentence data, which allows us to investigate to what extent the ALADIN approach can achieve its second goal: learning from sentence data without strong supervision such as word order and vocabulary.

In Section 2 we briefly describe the ALADIN system. In Section 3 we describe the dysarthric speech databases used for evaluation. In Section 4 we describe the experimental setup, and we discuss our results in Section 5. We present our conclusions and directions for future work in Section 6.

---

The research of J. F. Gemmeke was funded by the IWT-SBO project ALADIN (contract 100049). The STARDUST and VIVOCA projects were sponsored by the U.K. Department of Health New and Emerging Application of Technology (NEAT) programme. The views expressed in this publication are those of the authors and not necessarily those of the Department of Health. The research of S. Sehgal is currently supported by the EPSRC Programme Grant EP/I031022/1 (Natural Speech Technology).



**Fig. 1:** Visualisation of the non-negative matrix factorisation (NMF) approach to learn acoustic representations of semantic slot-values from sentences that are only weakly annotated at the utterance level.

## 2. ALADIN

### 2.1. Knowledge representation

Each spoken command, for example “turn on the television”, is associated with a possible action. A manual execution of the action would for example be pressing the standby button on the television remote control. Actions are represented using a *semantic frame* [12], a data structure that represents the semantic concepts that are relevant to the execution of the action and which end-users are likely to refer to in their spoken commands. Each semantic frame represents a possible action, and is composed of slots, which in turn contain slots or values. To continue the example, a semantic frame could contain two slots, `<device>` and `<action>`, allowing the values `<television, radio>` and `<on, off>`, respectively.

Internally, a semantic frame description is represented as a binary *label vector* indicating the presence or absence for all possible slot-values collected over all frames and slots. Using the example semantic frame, the command “turn on the television” would be represented as `[1 0 1 0]`.

### 2.2. Non-negative matrix factorisation

The ALADIN approach works by determining recurrent acoustic patterns in spoken commands, and is based on a non-negative matrix factorisation (NMF) approach [11, 10, 13]. NMF is a technique which decomposes a non-negative matrix into the product of two non-negative low-rank matrices [14, 15, 16, 17]. The approach is visualised in Fig. 1. First, the spoken command is converted into an utterance-based vector representation, the *acoustic representation*. In a nutshell, this representation is constructed for each utterance by making a histogram of the co-occurrences of Gaussian posteriors over time, with the Gaussian acoustic model obtained in advance. The acoustic model is estimated through unsupervised k-means clustering of the training data, followed by estimating a single full co-variance Gaussian on each cluster.

The collection of spoken training commands is concatenated into a matrix, the leftmost matrix in Fig 1, which is then factorised

by NMF into a matrix representing recurrent acoustic patterns (the *dictionary*), and a matrix of activations of these patterns over the training utterances. As visualised by the top half of the matrices in Fig 1, this factorisation is guided (regularized) by the label vectors to ensure that the obtained acoustic patterns correspond to slot-values within semantic frames. In addition to acoustic representations that are trained using supervision, a number of acoustic patterns are trained unsupervised to model acoustic phenomena occurring across sentences. These could include breathing sounds and silence, but also words for which no supervision is available, for example filler words.

### 2.3. Decoding

Decoding an observed utterance entails using NMF to find the combination of dictionary elements needed to represent the acoustic representation of the spoken command. Through the correspondence of these activations with the slot-values in semantic frames, we infer a semantic frame description of the observed utterance: for each slot whose cumulative slot-value activations exceeds a threshold, we assign the value with the largest activation.

## 3. SPEECH MATERIAL

In this work, we employ two datasets, VIVOCA (1 and 2) and STAR-DUST. The methods employed to collect this data are described in [5] and [3] respectively. All speakers had mild to moderate dysarthria. The speech was recorded directly onto either a laptop computer or a PDA hand-held computer.

### 3.1. VIVOCA

The vocabulary size, number of utterances and intelligibility assessment are shown in Table 1. The data from the VIVOCA project contains words that were used by the speakers to compose messages on voice output communication aid. The size of the vocabulary for each speaker varied according to the message building method

**Table 1:** Dysarthric speech databases used for evaluation. Intelligibility is denoted with E for less than 20% intelligibility, D for 20-50 % intelligibility and C for 50-90 % intelligibility. Starred labels are the result of informal listening tests, while non-starred labels are measured using the word-level intelligibility assessment procedure described in [5].

Speaker	VIVOCA													STARDUST		
	1	2	3	4	5	6	7	8	9	10	11	12	13	1	2	3
Vocabulary size	35	14	19	57	35	64	100	28	11	6	20	16	13	19	10	13
Total Utterances	1225	742	514	2956	1674	2821	4543	933	220	145	269	283	272	628	417	708
Intelligibility (%)	E*	D	E	E	E	E	E	E	C	E*	E	E*	D*	E	E	E

the speaker choose to use (see [5], section II B). For each speaker the message building method, and the input and output vocabularies were individually tailored to the needs and wishes of each participant. Generally, each word in the input vocabulary would map on to a short phrase. Longer phrases could be built up using combinations of words, meaning each allow sequence of words would produce a unique output sequence (or command).

### 3.2. STARDUST

The second and third datasets are based on data collected in the STARDUST project [3]. The second dataset is an isolated word recognition task using the same (`sil $word sil`) grammar as the VIVOCA data. It consists of three speakers and is constructed from the available training and adaptation data. The vocabulary size, number of utterances and intelligibility assessment are shown in Table 1.

The third dataset entails command & control sentences. Since the employed databases contain only few, if any sentence recordings we artificially constructed sentences by concatenating the waveforms of isolated words following a speaker-specific grammar. These grammars, shown in Fig. 2, were constructed to closely resemble those used in the STARDUST project, albeit somewhat simplified to account for shortages of some (isolated) words. While not a replacement for the full acoustic variation in real spoken sentences (albeit dysarthric speech may exhibit more pauses between words than regular speech), the data does suffice to evaluate the effectiveness of ALADIN approach of learning without segmentation/word order information.

A Voice Activity Detection (VAD) algorithm [18] was used to remove the silence in the isolated word waveforms prior to concatenation, although some pre-,inter-,and post-word silence remains. Every isolated word from the second database was (at most) only used once in the construction of the third database. The sentences were randomly generated while maintaining an as even distribution of words and grammar rules as possible. With respect to the isolated words STARDUST database (c.f. Table 1), the vocabulary of speaker 1 changed from 19 to 17 words, and the utterance counts for speaker 1-3 are now 260,204 and 490, respectively.

## 4. EXPERIMENTAL SETUP

### 4.1. ASR frontend

The conventional ASR front-end, referred to as ASR in the experimental results, employs left-to-right HMMs with 7 states per word, which yielded slightly better results than the 9 (non-emitting) states employed in [5]. Lower state counts, down to 3 states per word, were explored as well, but those lead to only very small improvements with few training samples, at the cost of a large performance decrease with more data. The acoustic vectors were 12 Mel-frequency cepstral coefficients (MFCCs) derived from a 26-channel filterbank with a 25 ms analysis window and 15 ms frame-rate. The models

were trained using the HMM toolkit [19] with the Baum-Welch algorithm.

### 4.2. ALADIN

The ALADIN system also operates on MFCC features. The system employs a VAD [18] to remove silence frames after feature extraction, as a proxy for the silence model employed by the ASR frontend. The mid-level acoustic representation, unique to each speaker, consists of 100 full-covariance Gaussians, trained on all speech material available for that speaker. Experiments with smaller and larger (50 to 200 Gaussians) acoustic representations did not yield substantially different results.

For the isolated word experiments, the semantic frame descriptions entail a single frame per word (without slots and slot-values). As a result, in this setup every individual word in the vocabulary is modelled by a single acoustic representation. The semantic frame descriptions for the sentence data were modelled after the grammars in Fig. 2 and are shown in Fig. 3. Other parameter settings were taken the same as in [11]. Most notably, the number of acoustic patterns that are trained unsupervised is 20% of the number of words (isolated word data) or slot-values (sentence data).

### 4.3. Evaluation procedure

We use the cross-validation technique described in [11]. In short, we divide the data in multiple blocks, with the constraints that each slot-value should occur in each block, and that the distribution of slot-values over blocks is as equal as possible. We evaluate with an increasing number of blocks used as training data, with the remaining blocks used as test data. The number of blocks is dependent on the amount of speech material and ranges from 10 to 6. To improve the statistical significance, we repeat the procedure with five different assignments of blocks to train and test data (folds). Evaluation is done using an F-score measure at the slot-value level, aggregated over all five folds. For more details we refer the reader to [11]. Note that for the isolated words datasets, the use of a single frame per word means the F-score is equal to the word classification accuracy.

## 5. RESULTS AND DISCUSSION

### 5.1. Isolated word data

The results of the evaluation on the isolated words VIVOCA database are shown in Fig. 4. When comparing the ASR and the ALADIN system we can observe that the ALADIN system achieves much higher F-scores at the beginning of the learning curves, ranging from 5 to 40% absolute. This remains true even for speakers (e.g. 4, 6, 7) for which the beginning of the learning curve represents dozens of examples per word — for speakers with much speech material the cross-validation procedure resulted in relatively initial training

```

$device1 = tv | disc | radio;
$device2 = film;
$device3 = video;
$state   = on | standby;
$control1 = sound | channel;
$control2 = play | stop;
$dir     = up | down;
$num     = one | two | three | four | five;
$cancel  = no;

```

```

$cstate      = sil $device1 sil $state sil;
$cctrl1     = sil $control1 sil $dir sil;
$cctrl2     = sil $device2 sil $control2 sil;
$cctrl3     = sil $device3 sil $nums sil;
$ccancel    = sil $cancel sil;

```

( \$cstate | \$cctrl1 | \$cctrl2 | \$cctrl3 | \$ccancel )

(a) STARDUST speaker 1

**\$device** = tv | radio | lamp;  
**\$state** = on | standby;  
**\$control** = volume | channel;  
**\$dir** = up | down;  
**\$cancel** = no;

```

$cstate      = sil $device sil $state sil;
$scntrl1    = sil $control sil $dir sil;
$ccancel    = sil $cancel sil;

```

( \$cstate | \$cctrl1 | \$ccancel )

(b) STARDUST speaker 2

```
$device      = tv | disc;
$state       = on | standby;
$control1    = volume | channel;
$control2    = play | stop | forward | back;
$dir         = up | down;
$cancel      = bugged;
```

```

$cstate      = sil $device sil $state sil;
$scntrl1     = sil $scntrl1 sil $dir sil;
$scntrl2     = sil $scntrl2 sil;
$ccancel     = sil $ccancel sil;

```

( \$cstate | \$ccntrl1 | \$ccntrl2 | \$ccancel )

(c) STARDUST speaker 3

**Fig. 2:** Grammars and vocabulary for each of the three speakers in the STARDUST sentence dataset.

<i>Frame</i>	<i>Slot</i>	<i>Value</i>
<b>on.off</b>	<action> <device>	on, off 1-3
<b>control</b>	<action> <function>	up, down vol, chan
<b>film</b>	<action>	play, stop
<b>video</b>	<action>	1-5
<b>cancel</b>	-	-

(a) STARDUST speaker 1

<i>Frame</i>	<i>Slot</i>	<i>Value</i>
<b>on_off</b>	<action> <device>	on, off 1-3
<b>control</b>	<action> <function>	up, down vol, chan
<b>cancel</b>	-	-

(b) STARDUST speaker 2

<i>Frame</i>	<i>Slot</i>	<i>Value</i>
<b>on_off</b>	<action> <device>	on, off tv, disc
<b>control</b>	<action> <function>	up, down vol, chan
<b>disc</b>	<action>	1-4
<b>cancel</b>	-	-

(c) STARDUST speaker 3

**Fig. 3:** Semantic frame descriptions for each of the three speakers in the STARDUST sentence data. Note that the slot-values do not directly correspond to the vocabulary in the grammars in Fig. 2, as they only represent human-readable tags of semantic concepts.

blocks. At the end of the learning curve, the systems perform comparably for most speakers, with ASR and ALADIN outperforming each other on some. For the isolated words STARDUST dataset in Fig. 5, we observe the same trends.

For both the ASR and the ALADIN system, we observe large performance differences between speakers at the end of the learning curves. Since the vocabulary size differs between speakers, one must be careful with direct comparisons. That said, intelligibility does seem to affect performance: the vocabulary for speaker 6 is much smaller (64 words) than for speaker 7 (100), but the latter achieves higher accuracies. The speakers with the best intelligibility assessment (2, 9 and 13) are among the best performing, but several other speakers (such as 1 and 5) perform comparably. At the same time, some speakers, such as 6 and 12, do not exceed F-scores of 70-75% even with substantial amounts of training data (hundreds of examples per word).

On isolated word data, the NMF-based learning almost boils down to a (Kullback-Leibler divergence weighted) averaging of the co-occurrence acoustic representations for each word, with as only difference the presence of acoustic patterns that are trained unsupervised. Small pilot studies suggest, however, that for this dataset the presence of these unsupervised acoustic patterns has only a minor impact on the result. It is interesting then, that such a simple utterance-based representation performs so much better than the ASR system.

Unfortunately, due to the differences between the NMF and ASR approaches it is difficult to compare aspects such as the total number of parameters. Even at the intermediate acoustic level of NMF, a comparison of the number of Gaussians employed does not tell the whole story: For example, the NMF system employs a single

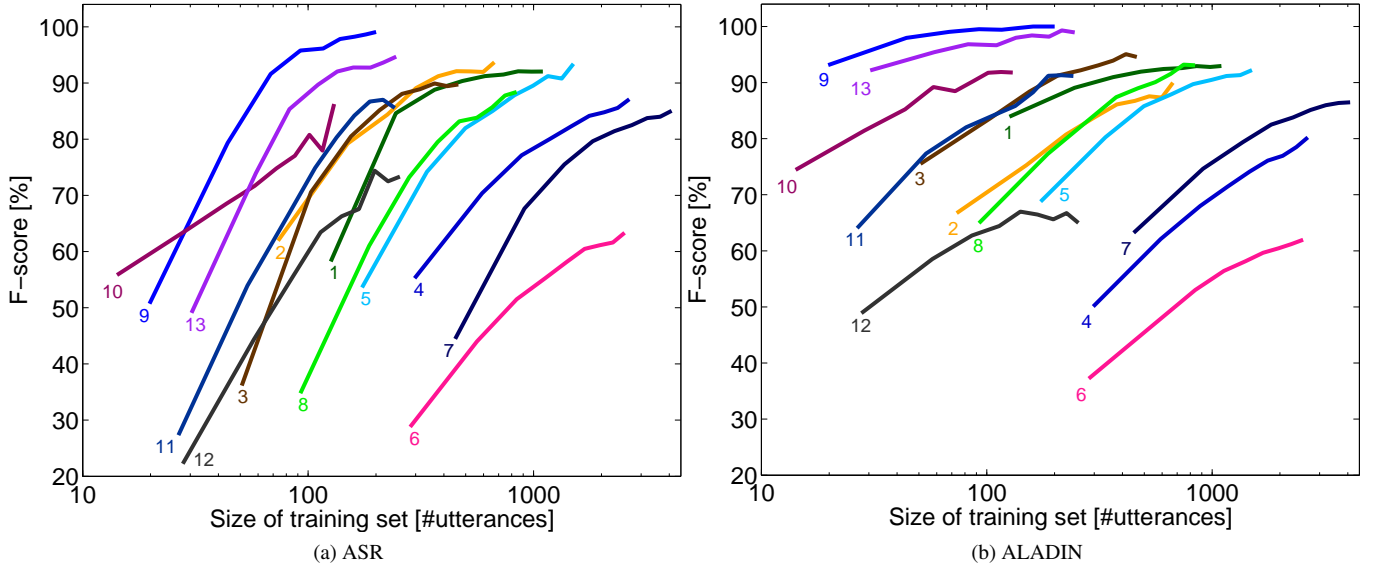
set of 100 full-covariance Gaussians, while the ASR system uses  $3 * 7 = 21$  Gaussians per word. Our experiments showed that in the NMF system, full-covariance Gaussians perform better than Gaussian mixtures, possibly due to the data-driven clustering approach used to train the Gaussians.

That said, experiments with more or less Gaussians (in the ALADIN system) and more or less states per word (in the ASR system) did not substantially improve results over the results presented here. This does suggest that utterance-based statistics may be a more stable word representation of highly variable dysarthric speech than HMM-based representations.

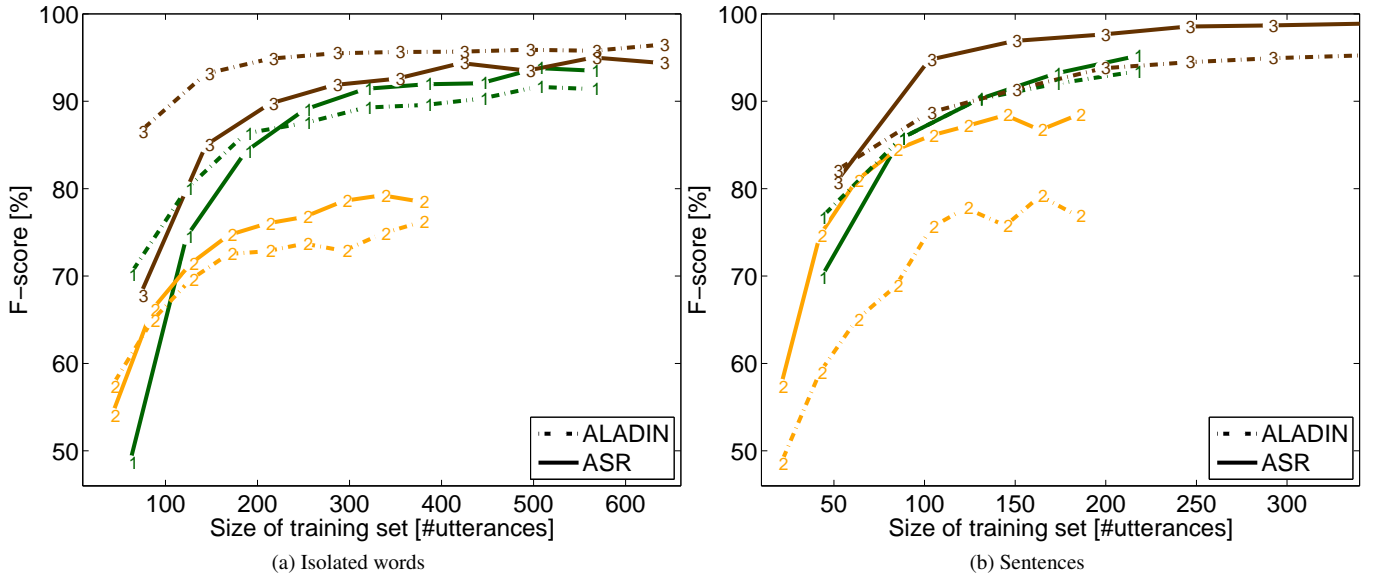
Although the usage scenario of the existing VIVOCA system is not the end goal of ALADIN, which encompasses not only fast learning but also learning from less detailed annotation, the ALADIN system may already be a viable approach to reduce the amount of training data needed. Naturally, if the end goal is only isolated word recognition, various alternative HMM-based approaches for learning with less data could also be explored, such as the use of tied Gaussians and subspace models.

## 5.2. Sentence data

On the STARDUST artificial sentence data displayed in Fig. 5b we observe a reverse of the isolated word results, with ASR now performing better than ALADIN even at the beginning of the learning curve for speakers 2 and 3, while performing comparably on speaker 1 at both ends. Direct comparison with the isolated word dataset is not possible, due to differences in the training size per cross-validation block, the vocabulary size (for speaker 1) and the recognition metric (for the sentence data the F-score is not equal to the word recognition



**Fig. 4:** VIVOCA isolated word recognition results per speaker as a function of the averaged number of utterances in the training set. The left panel displays the results obtained with the ASR system, a conventional GMM-HMM recognizer, whereas the right panel displays the results obtained with the NMF-based ALADIN framework. The graphs are displayed with a logarithmic horizontal axis to account for the large differences in the amount of training material. Numbers indicate the speaker index.



**Fig. 5:** STARDUST results per speaker as a function of the averaged number of utterances in the training set. The left panel displays the results obtained with isolated word recognition, whereas the right panel displays the results obtained with command & control sentences. Numbers indicate the speaker index.

accuracy).

That said, there are two aspects that contribute to the differences in performance of the ALADIN system and the ASR system between the isolated word data and the sentence data: 1) The ASR decoding results actually improve from the additional constraints imposed by the grammar, and 2) the ALADIN results decrease due to the difficulty of learning patterns from utterance-based representations for words that are never seen in isolation. Although more experiments

would be needed to isolate the contributions of these two aspects, it is encouraging that the ALADIN approach performs comparably to ASR for STARDUST speaker 1 even though that speaker has the most complex grammar. This confirms that when detailed annotation is not available, for example on under-resourced languages, or a scenario where the vocal interface is trained using only usage data (i.e., supervision consists of information such as the button press on a remote control), the ALADIN system may indeed be a viable approach

for vocal interfaces.

## 6. CONCLUSIONS

In this work we evaluate an approach to mine recurrent acoustic patterns from weakly supervised dysarthric speech data, to achieve two goals: 1) Achieving usable recognition accuracies with less training data, in order to minimize the initial effort of the target user, and 2) Achieving usable recognition accuracies with less detailed annotation - training a vocal interface using an unordered list of semantic concepts that are contained in the sentences, rather than a word-by-word transcript. Our contributions were the evaluation of both a speaker-dependent ASR approach and the ALADIN approach on large dysarthric speech databases, with speech of severely impaired speakers, and an evaluation of the performance on both isolated words and on C&C sentence data, which allows us to investigate to what extent the ALADIN approach can achieve its second goal: learning from sentence data without strong supervision such as word order and vocabulary.

The evaluations showed that on isolated word data the proposed approach achieves much higher accuracies when relatively few data is available, ranging from 5 to 40% absolute. With more data, the conventional ASR system and the proposed ALADIN approach perform comparably. We can conclude that for isolated words — the usage scenario of the existing VIVOCA system described in VIVOCA — the ALADIN system may already be a viable approach to reduce the amount of training data needed. In practical terms, this means users would be able to effectively use a VIVOCA system within days, rather than months of collecting speech material. For sentence data, more evaluation is needed, although it is impressive that the ALADIN approach performs comparably to ASR for STARDUST speaker 1 even though that speaker has the most complex grammar. Future work will focus on comparisons on sentence data from more speakers, real sentences, and with less constrained grammars and vocabulary.

## 7. REFERENCES

- [1] David Beukelman and Pat Mirenda, “Augmentative and alternative communication,” 2005.
- [2] J. Noyes and C. Frankish, “Speech recognition technology for individuals with disabilities,” *Augmentative and Alternative Communication*, vol. 8, no. 4, pp. 297–303, 1992.
- [3] M. S. Hawley, P. Enderby, P. Green, S. Cunningham, S. Brownsell, J. Carmichael, M. Parker, A. Hatzis, P. O’Neill, and R. Palmer, “A speech-controlled environmental control system for people with severe dysarthria,” *Medical Engineering & Physics*, vol. 5, no. 29, pp. 586 – 593, 2007.
- [4] Jort F. Gemmeke, B. Ons, N. Tessema, H. Van hamme, J. van de Loo, W. Daelemans, G. De Pauw, J. Huyghe, J. Derboven, L. Vugen, B. van Den Broeck, P. Karsmakers, and B. Vanrumste, “Self-taught assistive vocal interfaces : An overview of the ALADIN project,” in *Proc. INTERSPEECH*, 2013, pp. 1–5.
- [5] M.S. Hawley, S.P. Cunningham, P.D. Green, P. Enderby, R. Palmer, S. Sehgal, and P. O’Neill, “A voice-input voice-output communication aid for people with severe speech impairment,” *Neural Systems and Rehabilitation Engineering, IEEE Transactions on*, vol. 21, no. 1, pp. 23–31, Jan 2013.
- [6] H. Christensen, S. Cunningham, C. Fox, P. Green, and T. Hain, “A comparative study of adaptive, automatic recognition of disordered speech,” in *Proc Interspeech 2012*, Portland, Oregon, US, Sep 2012.
- [7] K. T. Mengistu and F. Rudzicz, “Comparing humans and automatic speech recognition systems in recognizing dysarthric speech,” in *Proceedings of the Canadian Conference on Artificial Intelligence*, 2011.
- [8] H. V. Sharma and M. Hasegawa-Johnson, “State transition interpolation and map adaptation for HMM-based dysarthric speech recognition,” in *HLT/NAACL Workshop on Speech and Language Processing for Assistive Technology (SLPAT)*, 2010, pp. 72–79.
- [9] F. Rudzicz, “Acoustic transformations to improve the intelligibility of dysarthric speech,” in *Proceedings of the Second Workshop on Speech and Language Processing for Assistive Technologies (SLPAT2011)*, 2011.
- [10] Jort F. Gemmeke, J. van de Loo, G. De Pauw, J. Driesen, H. Van hamme, and W. Daelemans, “A self-learning assistive vocal interface based on vocabulary learning and grammar induction,” in *Proc. INTERSPEECH*, 2012, pp. 1–4.
- [11] Bart Ons, Netsanet Tessema, Janneke van de Loo, Jort F. Gemmeke, Guy De Pauw, Walter Daelemans, and Hugo Van hamme, “A self learning vocal interface for speech-impaired users,” in *Proc. Workshop on Speech and Language Processing for Assistive Technologies (SLPAT)*, 2013, pp. 73–81.
- [12] Y. Wang and A. Acero, “Rapid development of spoken language understanding grammars,” *Speech Communication*, vol. 48, no. 3-4, pp. 390–416, 2006.
- [13] B. Ons, J. F. Gemmeke, and H. Van hamme, “Label noise robustness and learning speed in a self-learning vocal user interface,” in *Proc. of the International Workshop on Spoken Dialog Systems (IWSDS)*, Ermenonville, France, 2012.
- [14] D.D. Lee and H.S. Seung, “Learning the parts of objects by nonnegative matrix factorization,” *Nature*, vol. 401, pp. 788–791, 1999.
- [15] J. Eggert and E. Korner, “Sparse coding and NMF,” in *Neural Networks, 2004. Proceedings. 2004 IEEE International Joint Conference on*, 2004, vol. 4, pp. 2529–2533 vol.4.
- [16] Yu-Xiong Wang and Yu-Jin Zhang, “Nonnegative Matrix Factorization: A comprehensive review,” *Knowledge and Data Engineering, IEEE Transactions on*, vol. 25, no. 6, pp. 1336–1353, 2013.
- [17] Hyekyoung Lee, Jiho Yoo, and Seungjin Choi, “Semi-supervised nonnegative matrix factorization,” *Signal Processing Letters, IEEE*, vol. 17, no. 1, pp. 4–7, 2010.
- [18] Javier Ramírez, Juan Manuel Górriz, José Carlos Segura, Carlos G. Puntonet, and Antonio J. Rubio, “Speech/non-speech discrimination based on contextual information integrated bispectrum LRT,” *Signal Processing Letters, IEEE*, vol. 13, no. 8, pp. 497–500, 2006.
- [19] Steve Young, Gunnar Evermann, Dan Kershaw, Gareth Moore, Julian Odell, Dave Ollason, Valtcho Valtchev, and Phil Woodland, “The htk book,” *Cambridge University Engineering Department*, vol. 3, 2002.